

Technocolabs Data Science Internship

DA Project Report

TITLE: Lending Club's Loan Repayment Prediction

AIM:

The main purpose of this project is to develop machine learning models (Logistic Regression and Neural Networks) on Lending Club's data and deploy them for wider general-purpose usage.

ABSTRACT:

Traditionally, loan-level risk is measured as credit risk—the probability of default to measure the expected loss. Using machine learning techniques, we modeled credit risk and expected payoff maximization on the ROC, to help Lending Club optimize their risk. The models used here are Logistic Regression and Neural Networks.

This project aims to analyze the loans that were paid off in full or charged off.

INTRODUCTION:

Lending Club's dataset contains extensive information of their customer's loan status and many other metrics which help in the analysis of its data. The goal here is to build a model using this data and use that model to accurately predict whether the loans were either fully paid or charged off.

OVERVIEW:

1. Examine and Clean the Dataset
2. Visualize the dataset for gaining clarity
3. Build and Train the model
4. Visualize the Model's end Report

DATASET:

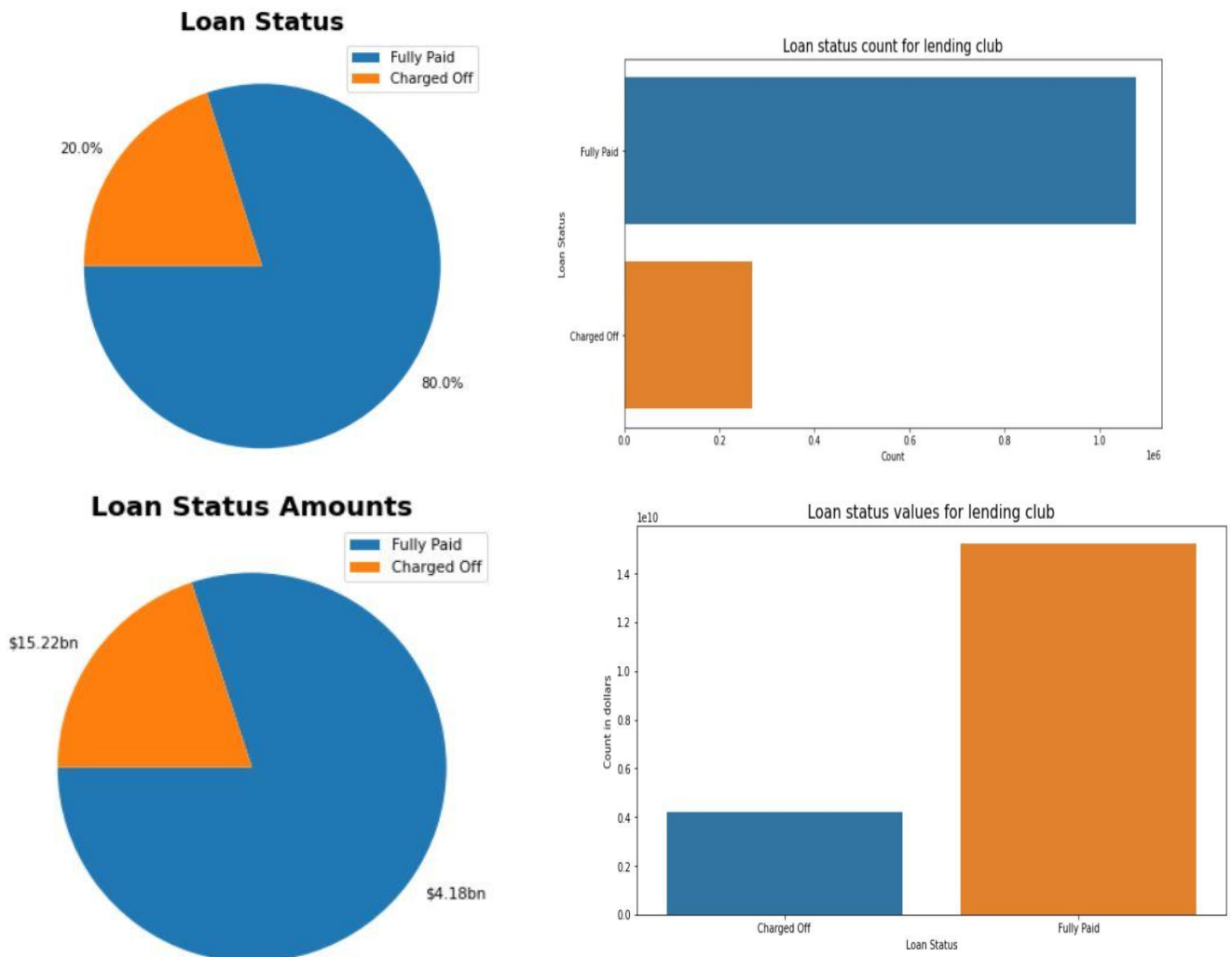
The Lending Club Dataset Contains 2260701 rows by default, which detail the number of customers and it has 151 columns which contain all details with respect to their lending club loan data requirements.

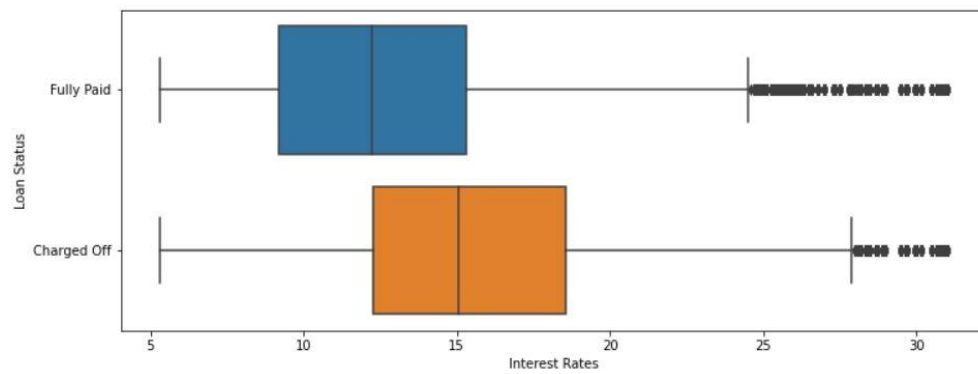
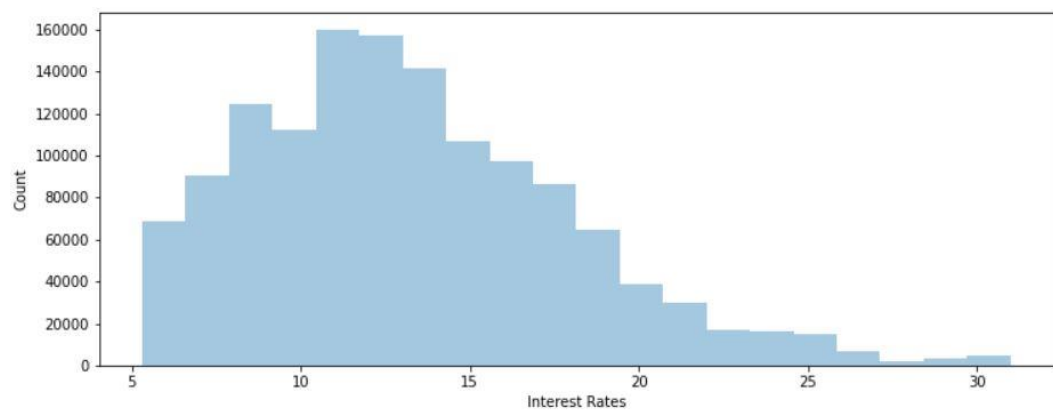
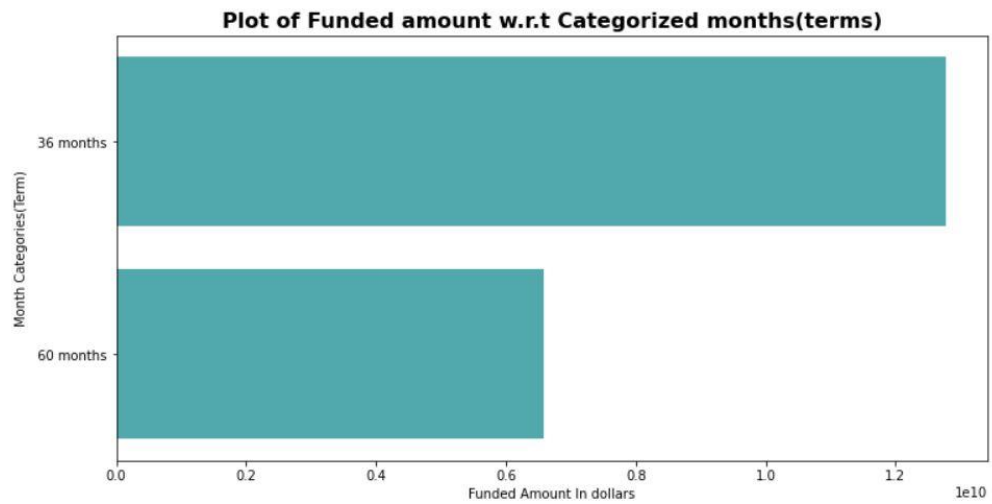
1) EXAMINE AND CLEAN THE DATASET:

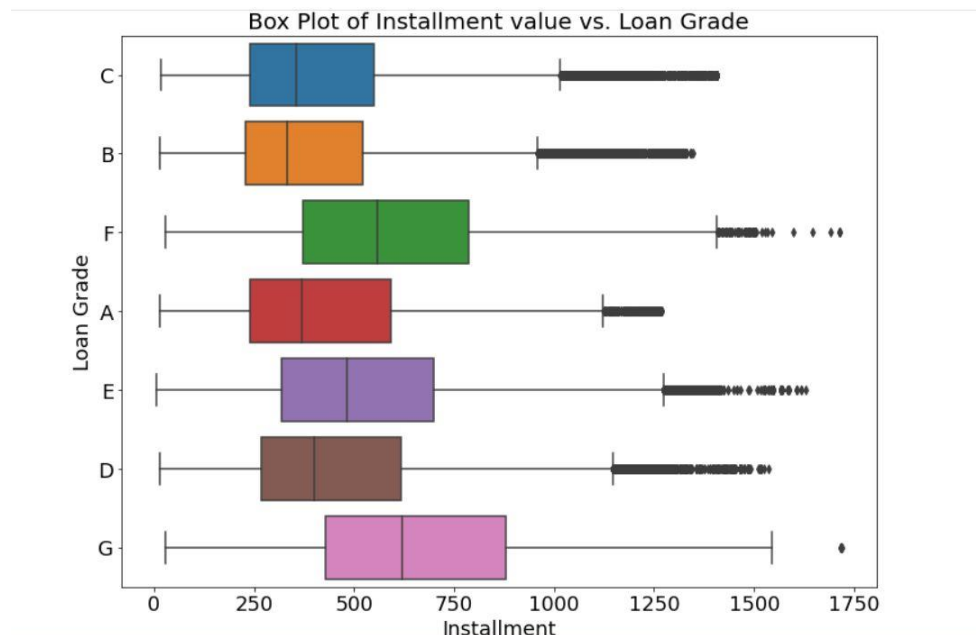
- The dataset is readily available from the bank's own website and via Kaggle.com as well. The data is in the form of .csv(comma-separated values) file. These files need to be loaded in the folder which contains the ipython notebook which will be used for analysis and model building.
- The dataset is then analysed for any discrepancies(mostly null values) and the percentage of null values in every column must be observed. If null values in any column exceeds 70%, then that column needs to be dropped to maintain accuracy.
- Then to analyze the loan status of the dataset, we label the values "Charged Off" as 0 and "Fully Paid" as 1 for proper classification

2) VISUALIZING THE DATASET:

- Here we visualize the dataset parameters using matplotlib and seaborn modules
- The charts produced are:

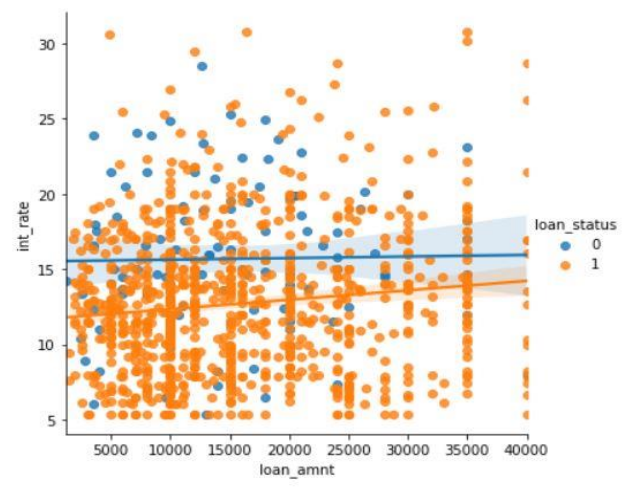
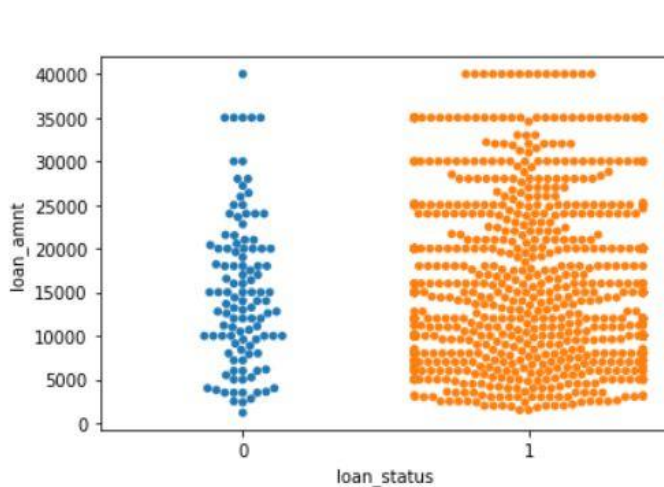






The below graphs are:

- Swarmplot of loan amount v/s loan status
- Lmplot of interest rate v/s loan amount

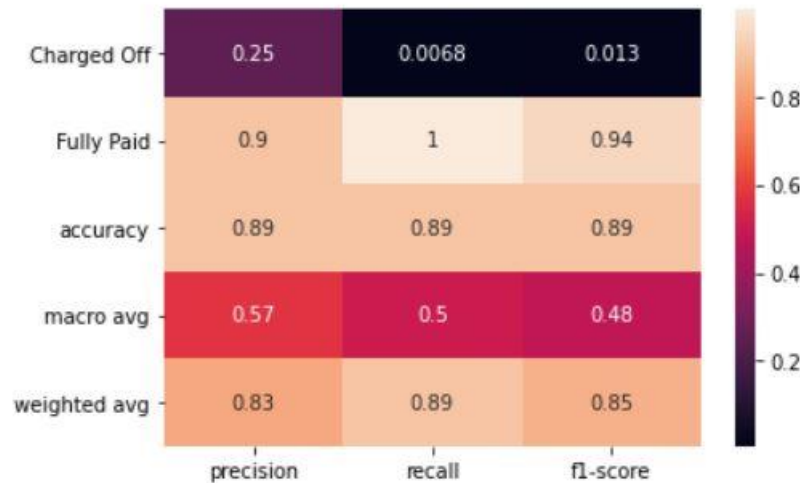


3) BUILDING AND TRAINING THE MODEL:

- The logistic regression model needs to be built using numeric data only. Hence the first step is to encode all the columns with non-numeric data types (they are mostly object data type) and assign numeric labels. This is done using One hot encoding.
- Next we assign the target values (y) to be the labelled column of loan status, wherein charged off is labelled as 0 and fully paid is labelled as 1. The feature matrix (X) consists of all the encoded columns which will be used for training the model.
- In the next step we split the data into training and testing sets using the `train_test_split` function.
- After obtaining the train and test sets we fit the data into the logistic regression model and let it train. After the training ends, we use the model to test it on the test set and obtain predicted values.

4) VISUALIZE THE MODEL'S END REPORT:

- To obtain the accuracy and other metrics of the model, we print out its classification report and subsequently plot a heatmap for easily visualizing it.
- The heatmap is as follows:



- Here as we can see, the model is approximately 89% accurate on unseen data.

DEPLOYMENT OF THE TRAINED AND TESTED MODEL:

- Deployment can be carried out by both flask and streamlit applications. Here we used streamlit.
- In Streamlit, first we saved our model using in a pickle file format(.pkl) and loaded it into a new python file. There we specified html parameters for the look of the end application and created a function which used this model to display the predictions. The predictions would be displayed as either “Charged Off” or “Fully Paid”.
- Then we finally deployed our streamlit model using Heroku for open world use.

Lending Club Loan Prediction ML App

Loan Amount

4000000.00

Interest Rate

17.00

Annual Income

800000.00

Debt To Income Ratio

6.00

FICO Range(Low)

699.99

Number of derogatory public records

5.00

Revolving Line Utilization Rate

5.98

Total Number of Credit lines

2.97

Months Since Most Recent Installment Accounts Opened

12.00

Number of Mortgage Accounts

0.00

Months Since Most Recent Bankcard Account Opened

2.00

Months Since Most Recent Inquiry

0.00

Predict

Your Loan is Charged Off